

# Statistische Entscheidungs Theorie

Dozent: Dr. Zoran Nikolić

Machine Learning Seminar



# Gliederung

- Machine Learning Einführung
- K-Nächste-Nachbarn
- Least Squares
- Vergleich der Modelle
- Fazit

# Machine Learning

- Ziel: Finde Zusammenhang zwischen Input Daten  $X$  und Output Daten  $Y$  um möglichst genaue Vorhersagen zu treffen:

- Wir suchen: 
$$\hat{Y} = \hat{f}(X) \approx f(X) = Y$$

- Input Daten  $X \in \mathbb{R}^{N \times K}$  (bzw. Zufallsvariable) werden auch **Features** oder **unabhängige Variablen** genannt
- $N \hat{=}$  #**Beobachtungen**
- $K \hat{=}$  #**Merkmale**
- Output Daten  $Y \in \mathbb{R}^N$  (bzw. Zufallsvariable) werden auch **Targets** oder **abhängige Variablen** genannt

# Vorhersagen

- **Quantitative** Vorhersagen: Output Daten  $Y$  liegen in einem stetigen Wertebereich
  - Bsp.: Vorhersage von Aktienkurswerten oder Temperatur
  - Bsp.:  $Y(X) = 2X + 5$
  - Wird auch als **Regression** bezeichnet

# Vorhersagen

- **Qualitative** Vorhersagen: Output Daten  $Y$  bestehen aus Klassen
  - Bsp.: Vorhersage ob Aktienkurs steigt oder sinkt oder ob Patient stirbt oder überlebt
  - Für 2 Klassen:  $Y(X) = \begin{cases} 1 & \text{falls } X \in \text{Klasse 1} \\ 0 & \text{falls } X \in \text{Klasse 0} \end{cases}$
  - Für  $K > 2$  Klassen :  $Y \in \mathbb{R}^{N \times K}$  mit  $K$  binären Vektoren
  - Wird auch als **Klassifikation** bezeichnet

# Vorhersagen

- Genaue Vorhersagen  $\Leftrightarrow$  Kleiner Fehler zwischen vorhergesagten und echten Werten  $Y - \hat{Y}$

- **Expected Prediction Error (EPE):**

$$EPE(f) = E \left( (Y - f(X))^2 \right)$$

- Es gilt: Minimum der Funktion  $EPE(f)$  liegt bei

$$f(X) = E(Y|X)$$

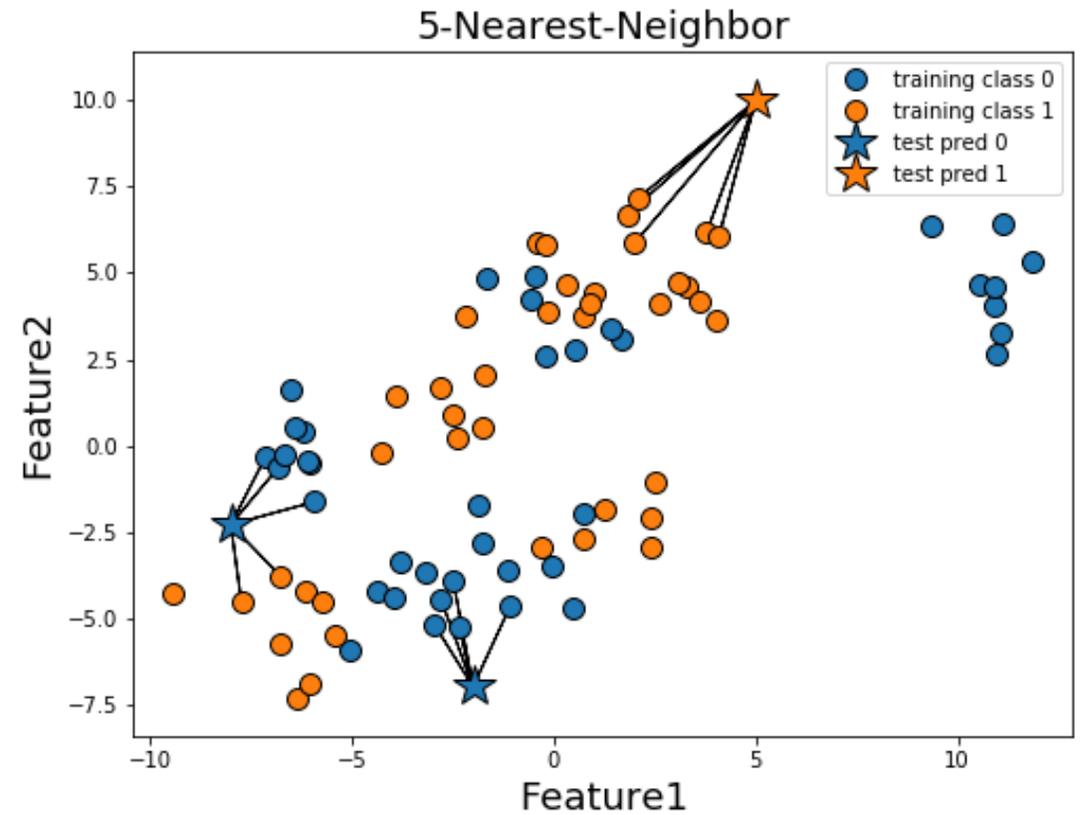
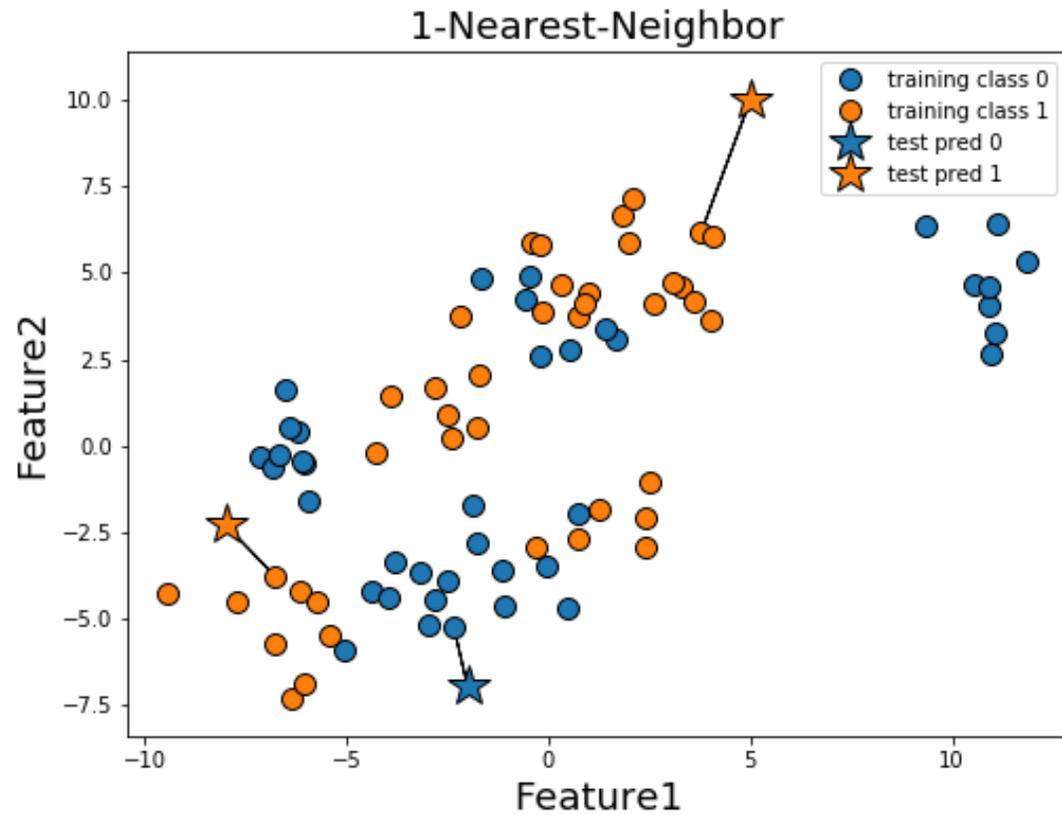
# K-Nächste-Nachbarn

- Ersetze Erwartungswert durch lokalen Mittelwert:

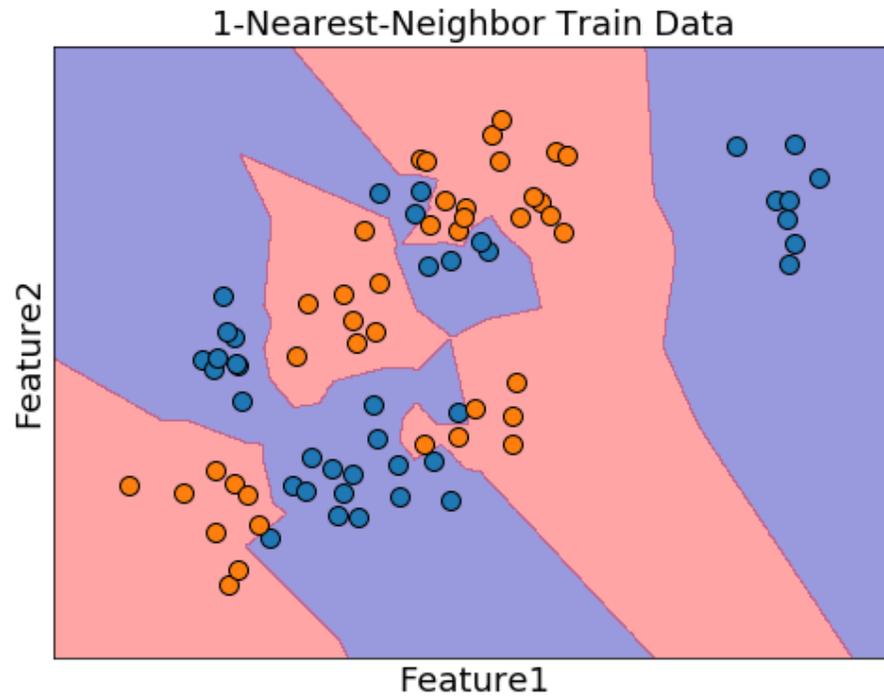
$$f(X) = E(Y|X) \approx \frac{1}{k} \sum_{x_i \in N_k(x)} y(x_i) =: \hat{f}(x)$$

$N_k(x) \hat{=}$  Menge der  $k \in \mathbb{N}$  nächsten Nachbarn, bzgl. euklidischer Norm

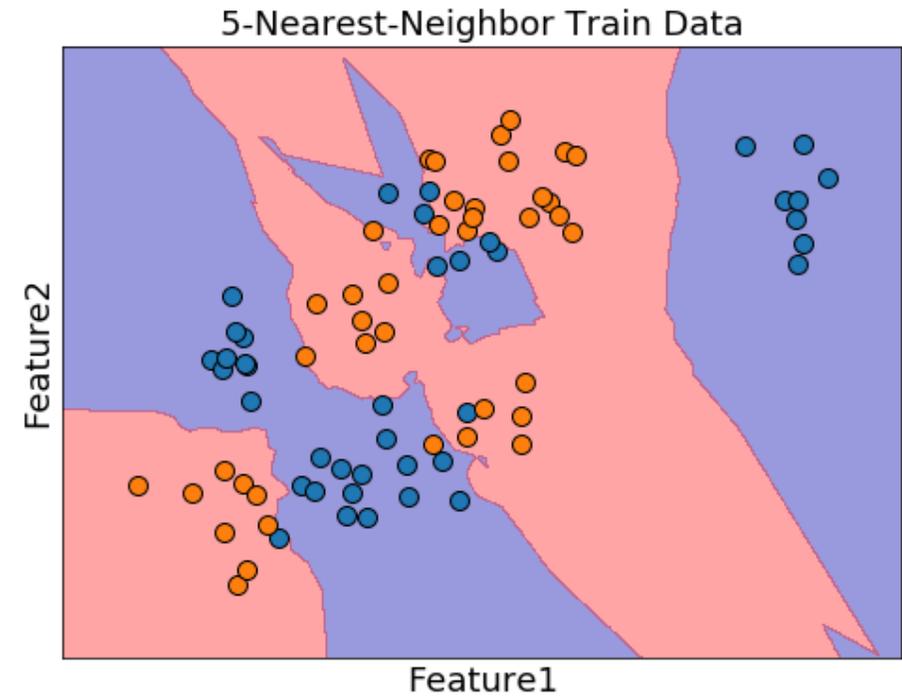
# K-Nächste-Nachbarn



# K-Nächste-Nachbarn

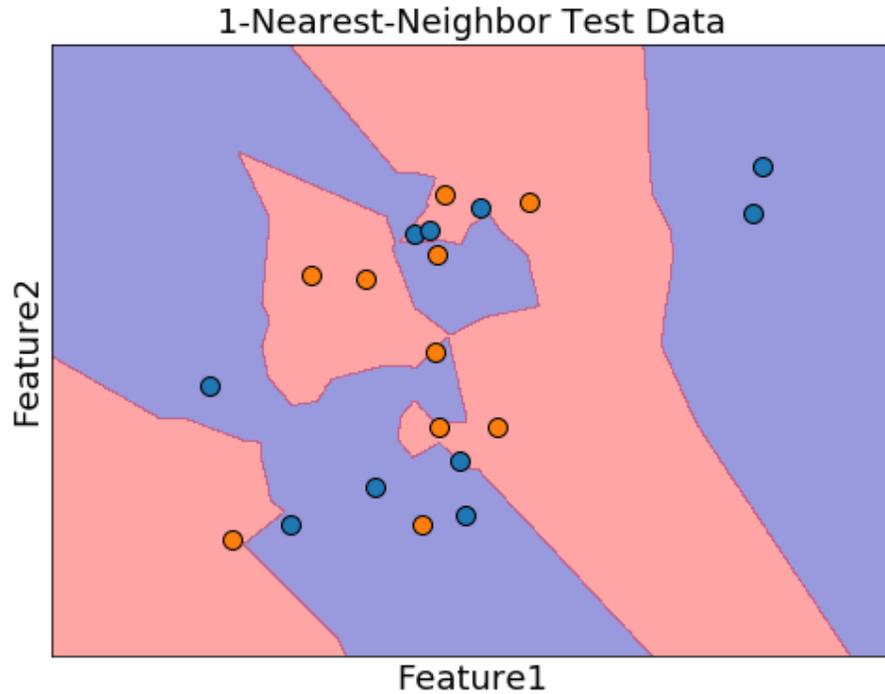


Genauigkeit: 100%

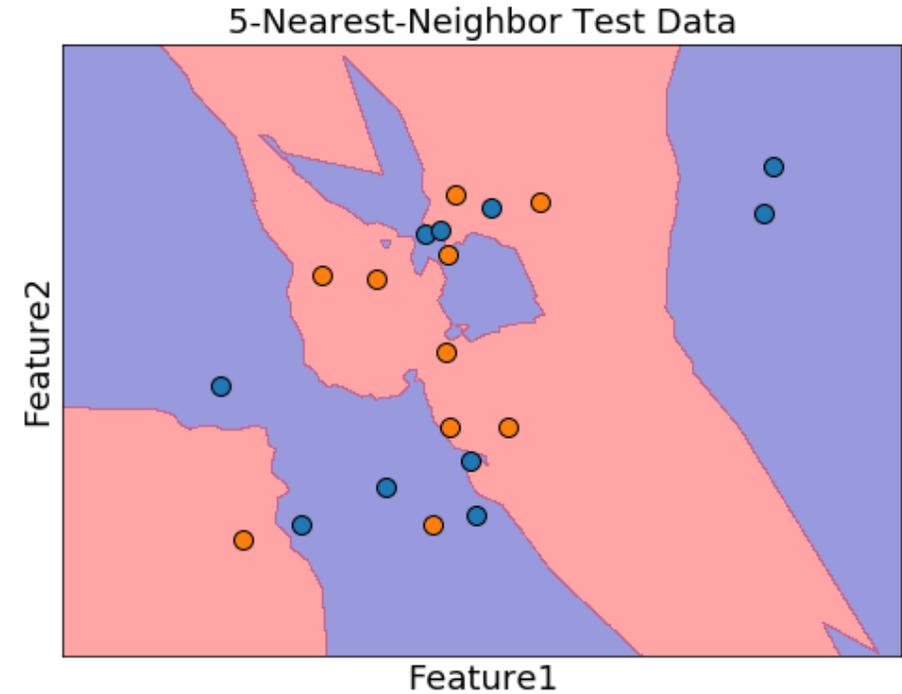


Genauigkeit: 94%

# K-Nächste-Nachbarn - Testdaten

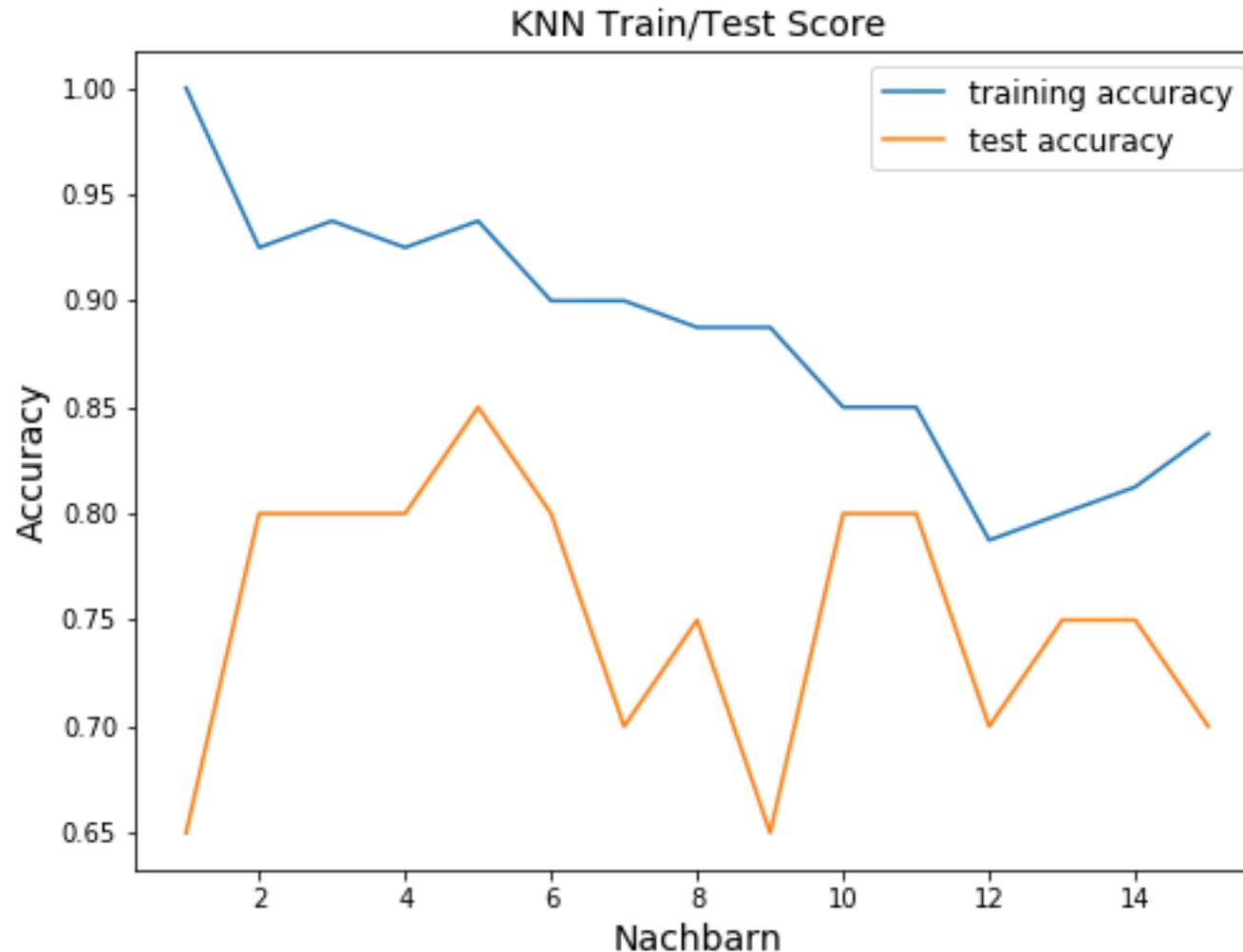


Genauigkeit: 65%



Genauigkeit: 85%

# K-Nächste-Nachbarn – Train/Test Daten



- $k = 1$ : zu stark an den Trainingsdaten angepasst
- $k = 12$ : Modell nicht Komplex genug

# Least Squares

- Annahme: Linearer Zusammenhang

$$Y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\Rightarrow f(X) = E(Y|X) = E(X\beta + \varepsilon|X) = X\beta + E(\varepsilon|X) = X\beta$$

- Einsetzen in  $EPE(f(X))$  und Erwartungswert durch Mittelwert ersetzen:

$$EPE(f(X)) = E((Y - X\beta)^2) \approx \frac{1}{N} \sum_{i=1}^N (y_i - x_i^T \hat{\beta})^2$$

# Least Squares

- Alternative Matrixschreibweise:

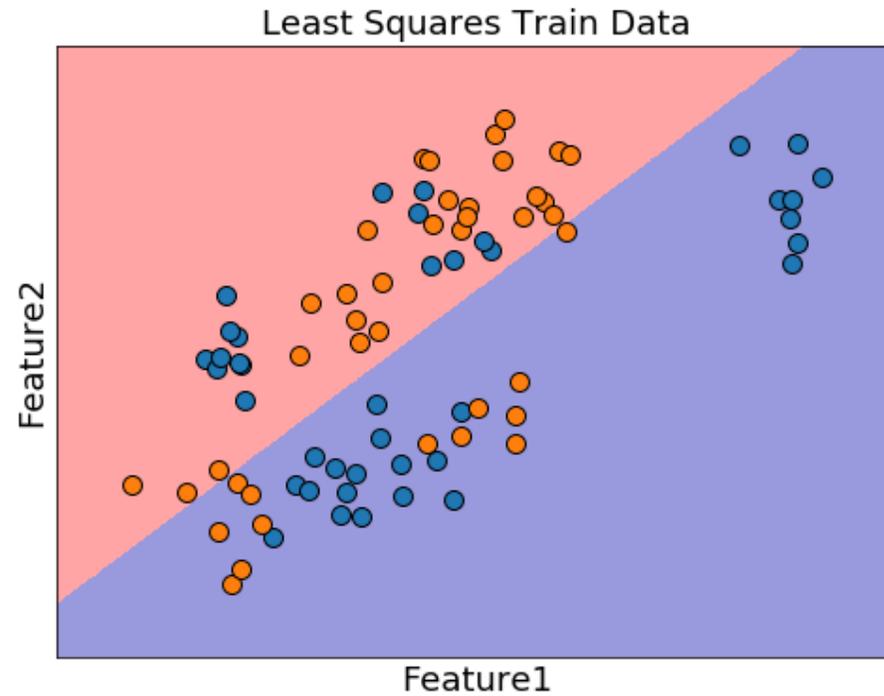
$$EPE(X) \approx \frac{1}{N} \sum_{i=1}^N (y_i - x_i^T \hat{\beta})^2 = \frac{1}{N} (y - X\hat{\beta})^T (y - X\hat{\beta})$$

- Ableitung nach  $\beta$  gleich 0 setzen und umformen ergibt:

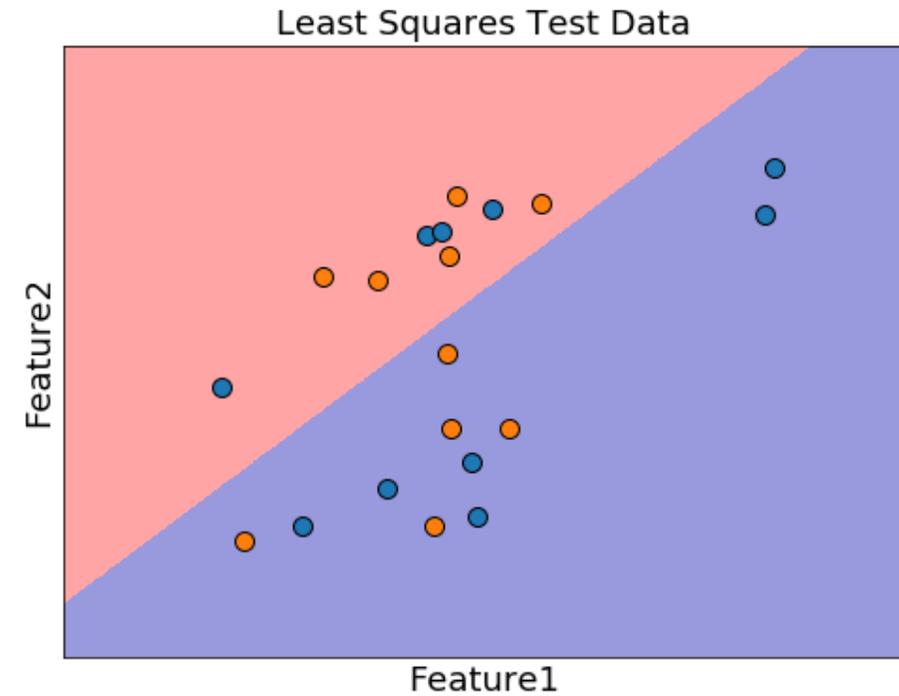
$$\frac{\partial EPE(X)}{\partial \hat{\beta}} = 0 \quad \Leftrightarrow \quad \hat{\beta} = (X^T X)^{-1} X^T y$$

$$\Rightarrow \hat{y}_i = \begin{cases} 1 & \text{falls } x_i^T \hat{\beta} \geq 0.5 \\ 0 & \text{falls } x_i^T \hat{\beta} < 0.5 \end{cases}$$

# Least Squares – Klassifikation

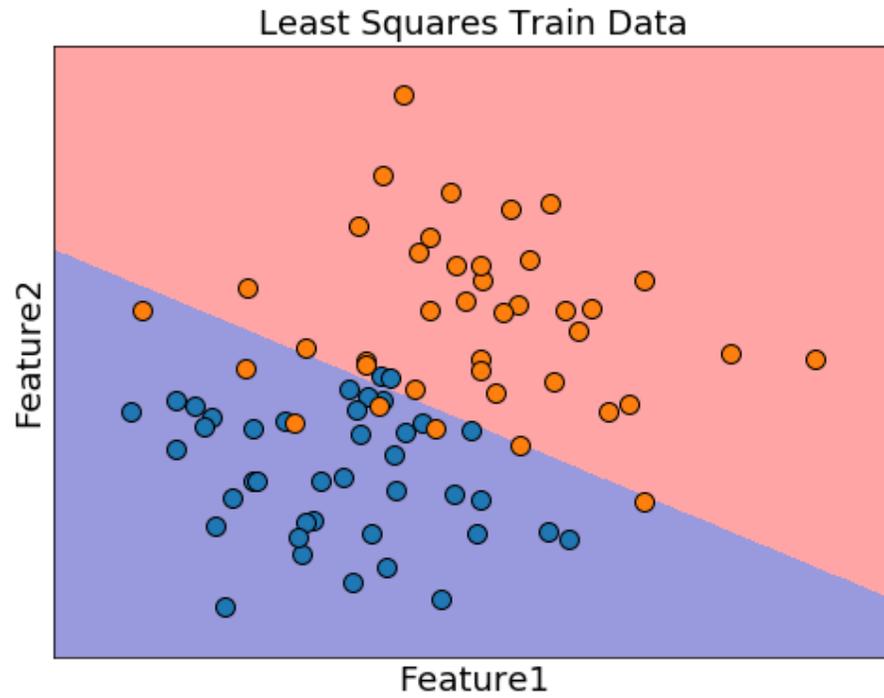


Genauigkeit: 63.75%

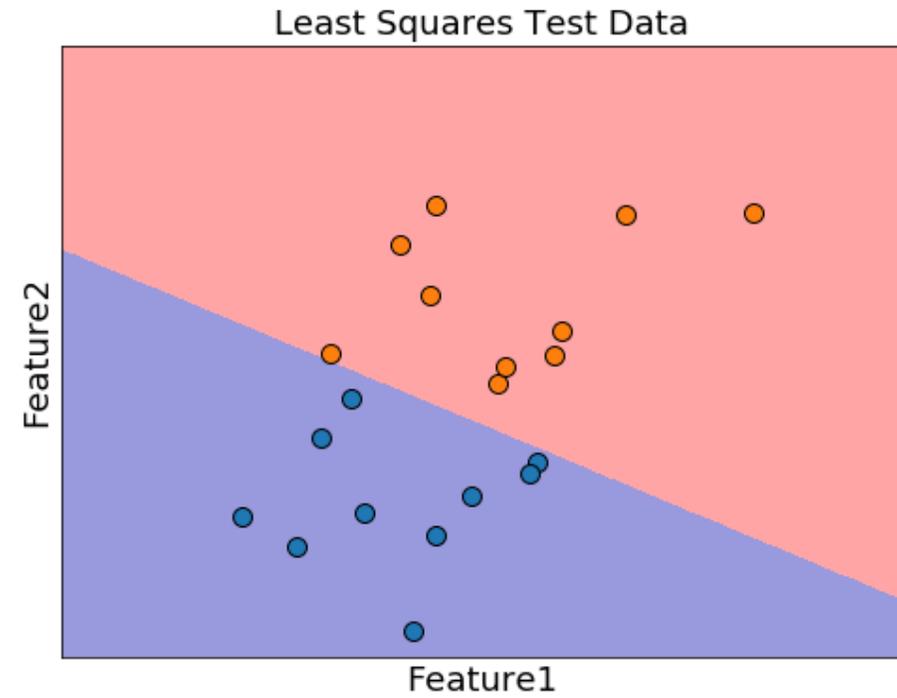


Genauigkeit: 55%

# Least Squares – Bessere Verteilung

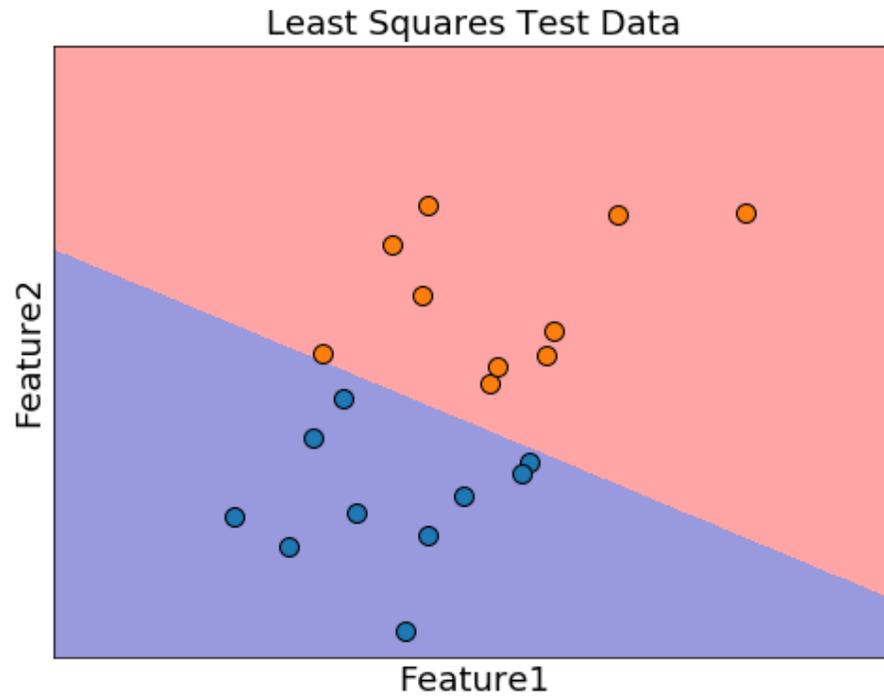


Genauigkeit: 88.75%

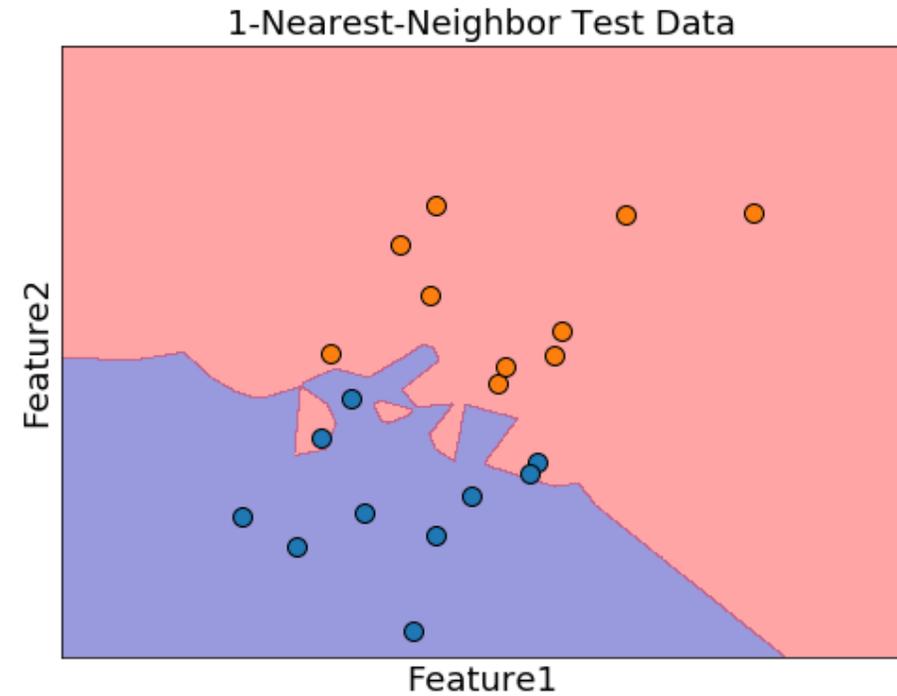


Genauigkeit: 100%

# LS vs. KNN – Klassifikation



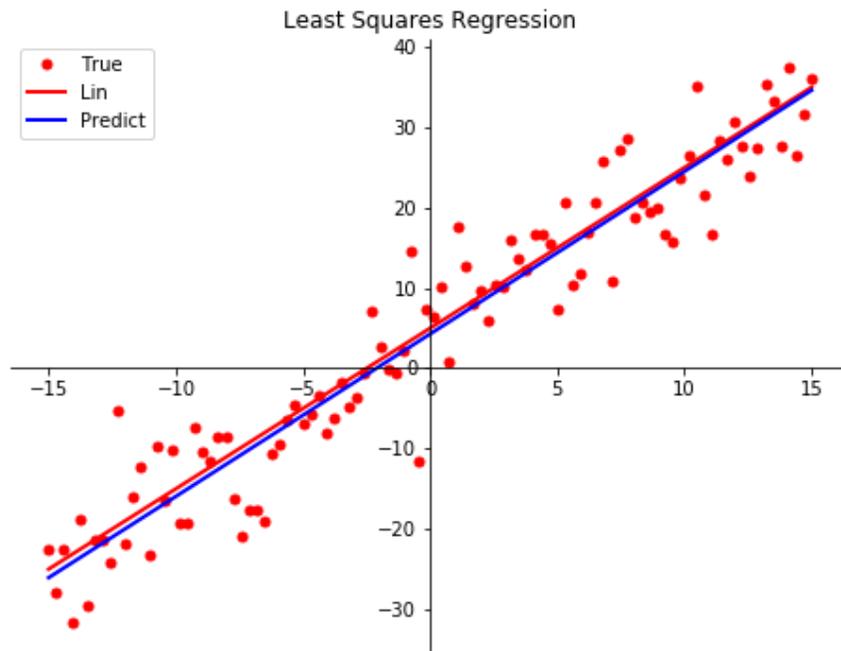
Genauigkeit: 100%



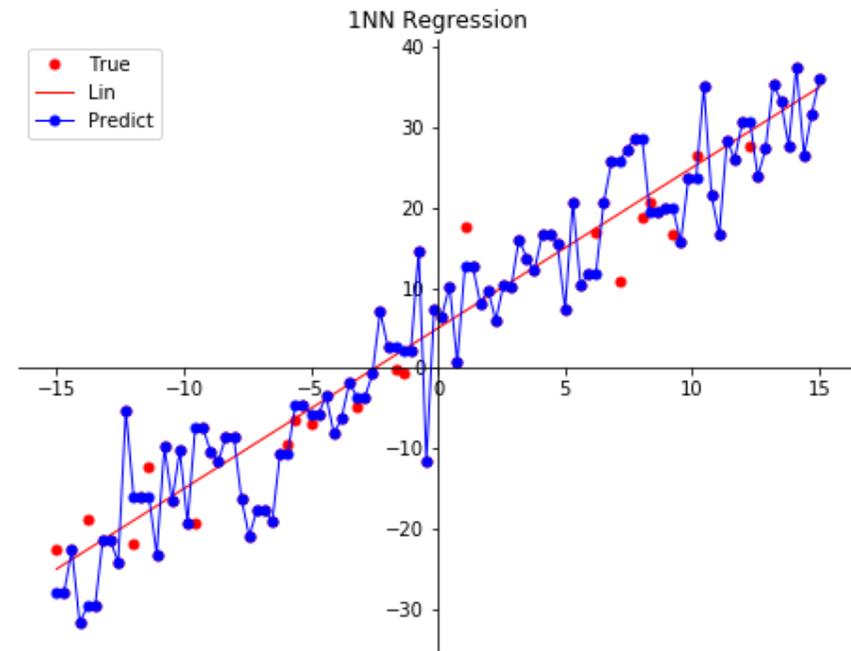
Genauigkeit: 85%

# LS vs. KNN – Regression

$$y = 5 + 2x + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0,25)$$



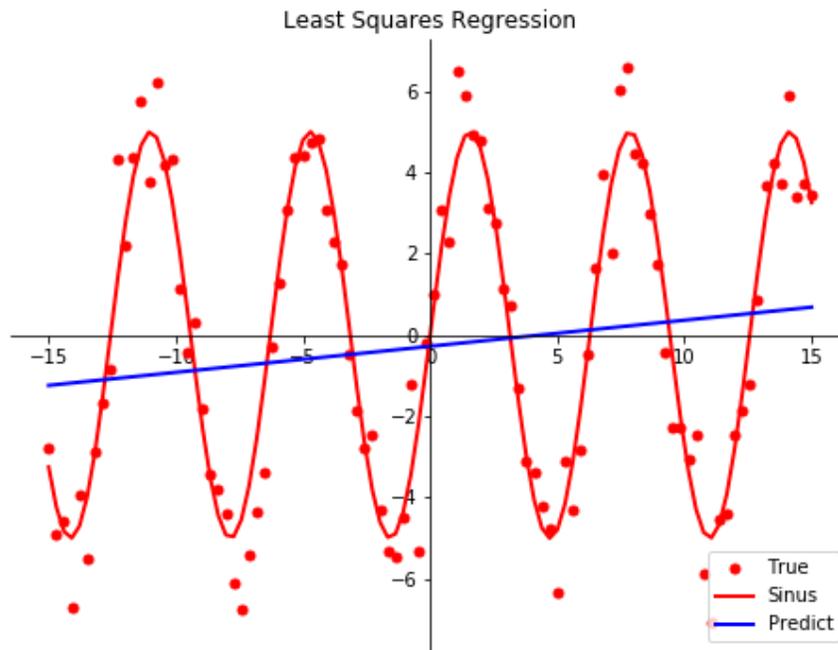
Mean Squarred Error: 18.4



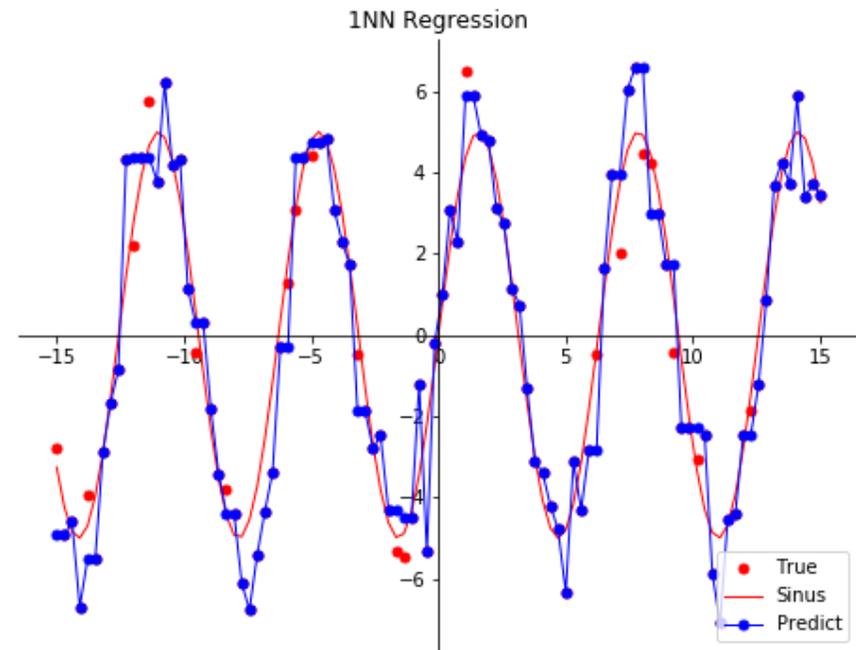
Mean Squarred Error: 37.7

# LS vs. KNN – Regression

$$y = \sin(x) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0,1)$$



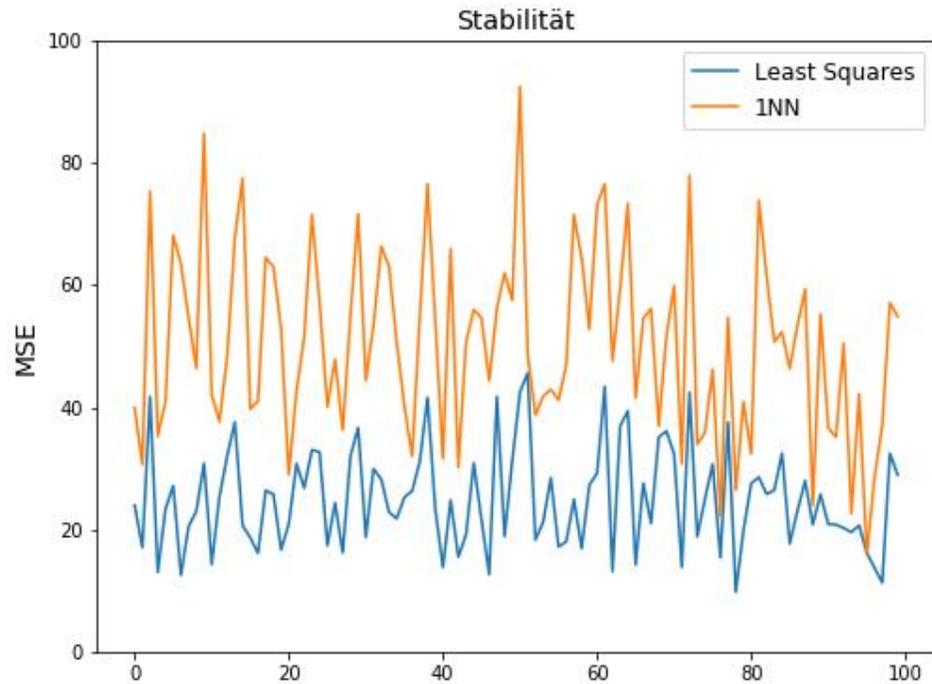
Mean Squared Error: 13.4



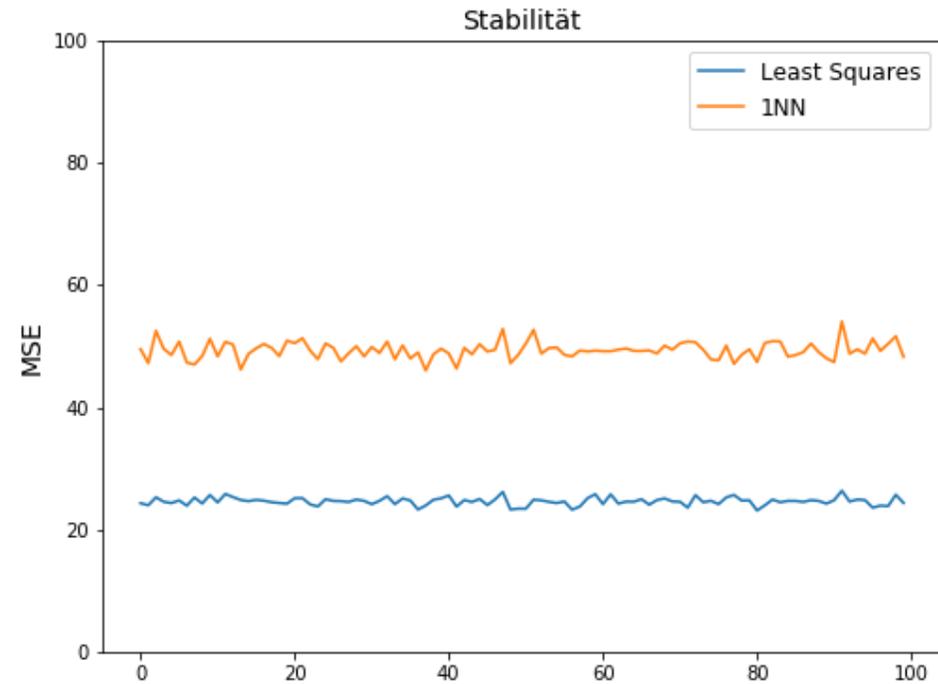
Mean Squared Error: 2.2

# Stabilität

$$y = 5 + 2x + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 25)$$



$N = 100$



$N = 10000$

X-Achse: Unterschiedliche Aufteilungen der Daten in Trainings- und Testdaten

# Bias-Varianz Dekomposition

- Mean Squared Error an der Stelle  $x_0$  auswerten für:

$$y = f(x) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 25)$$

$$\begin{aligned} \text{MSE}(x_0) &= E_{\mathcal{T}}((y_0 - \hat{y}_0)^2) \\ &= (f(x_0) - E_{\mathcal{T}}(\hat{y}_0))^2 + (E_{\mathcal{T}}(\hat{y}_0^2) - E_{\mathcal{T}}(\hat{y}_0)^2) + \sigma^2 \\ &= \text{Bias}_{\mathcal{T}}^2(\hat{y}_0) + \text{Var}_{\mathcal{T}}(\hat{y}_0) + \sigma^2 \end{aligned}$$

# Bias-Varianz Dekomposition

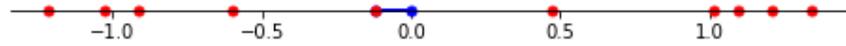
- $Bias_{\mathcal{T}}(\hat{y}_0)$ : Misst wie gut das Modell den deterministischen Teil von  $y_0$  approximiert
- $Var_{\mathcal{T}}(\hat{y}_0)$ : Varianz der unterschiedlichen Trainings/Test Aufteilung. Wird kleiner je größer  $N$  ist.
- $\sigma^2$ : Varianz des Fehlerterms  $\varepsilon$

# Bias-Varianz Dekomposition

- Least Squares:
  - $Var_{\mathcal{T}}(\hat{y}_0)$  wird klein für große  $N$
  - $Bias(\hat{y}_0)$  ist klein, wenn alle Annahmen erfüllt sind, ansonsten ist auch Least Squares verzerrt.
  
- KNN:
  - $Var_{\mathcal{T}}(\hat{y}_0)$  wird klein für große  $N$ , jedoch größer als bei Least Squares
  - $Bias(\hat{y}_0)$  wird groß, falls der nächste Nachbar weit entfernt ist, besonders bei hohen Dimensionen

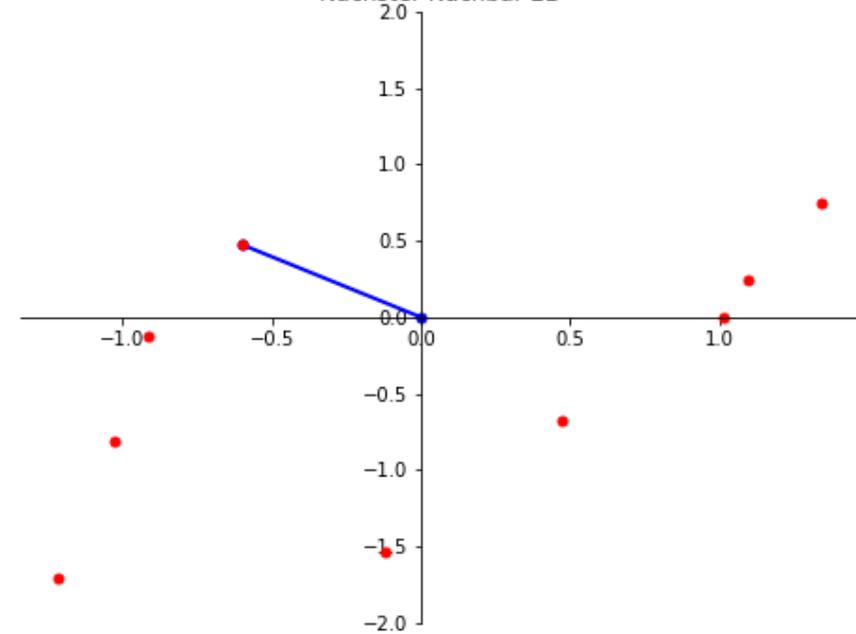
# KNN 1D vs. 2D

Nächster Nachbar 1D



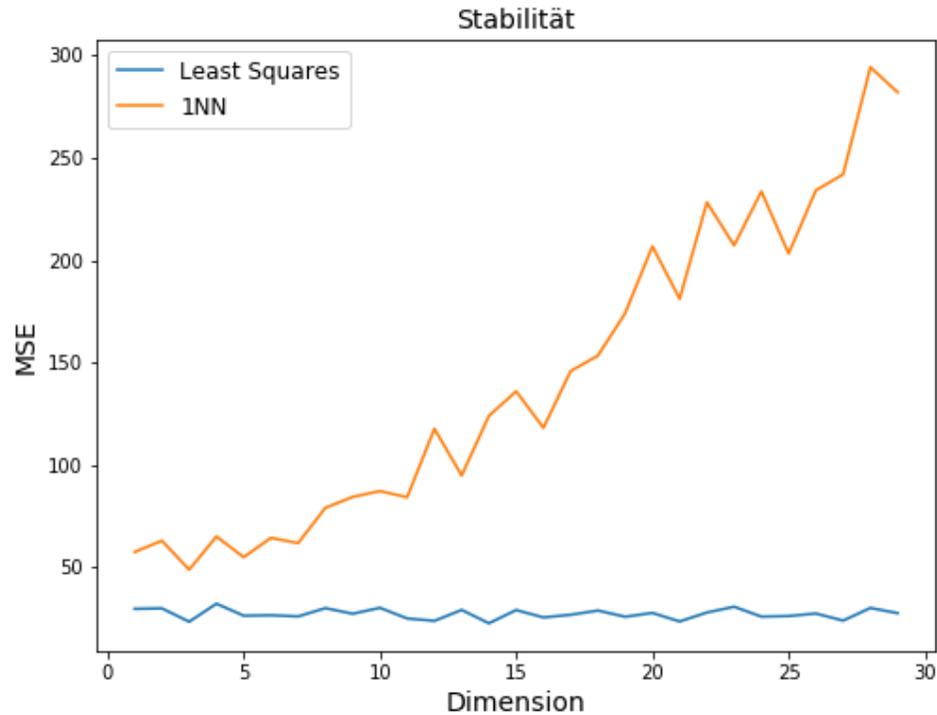
$$||x_0 - x_1|| = 0.12$$

Nächster Nachbar 2D

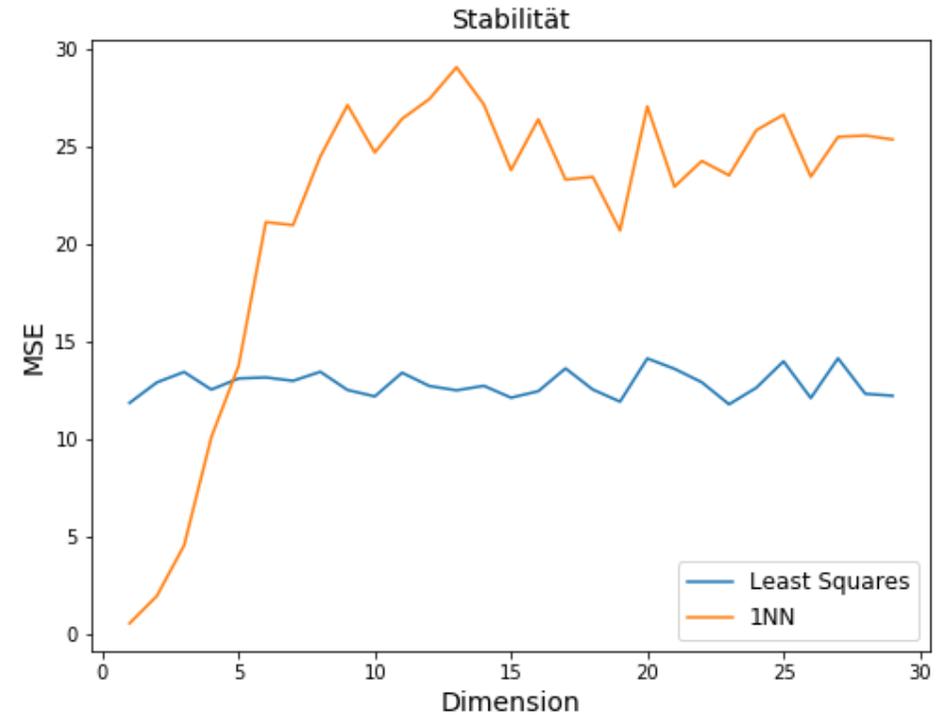


$$||x_0 - x_1|| = 0.76$$

# LS vs. KNN – Hohe Dimensionen



$$f(x) = 5 + \sum_{i=1}^d 2x_i$$



$$f(x) = 5\sin(\sum_{i=1}^d 2x_i)$$

# Least Squares – Vor- und Nachteile

- Vorteile:
  - Stabil und unverzerrt
  - Gilt auch für höhere Dimensionen
- Nachteile:
  - Erfolg ist an oft unrealistischen Annahmen gebunden
  - Nur einfache Zusammenhänge können erklärt werden

# KNN – Vor- und Nachteile

- Vorteile:
  - Keine Annahmen an die Daten notwendig
  - Komplexe Zusammenhänge können erfasst werden
- Nachteile:
  - Starke Verzerrung und Instabil
  - Vor allem für hohe Dimensionen

# Fazit

- Wahl des Modells hängt stark von den zugrundeliegenden Daten ab
- Untersuchung der Daten ist ein wichtiger Bestandteil des Machine Learning Bereichs
- Anzahl der Beobachtungen sollte möglichst groß sein
- Anzahl der unterschiedlichen Features sollte möglichst reduziert werden

Vielen Dank für Ihre Aufmerksamkeit.

